



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

The grant is dead, long live the data - migration as a pragmatic exit strategy for research data preservation [version 1; peer review: 2 approved].

Citation for published version:

Zielinski, T, Hay, J & Millar, AJ 2019, 'The grant is dead, long live the data - migration as a pragmatic exit strategy for research data preservation [version 1; peer review: 2 approved].', *Wellcome Open Research*, vol. 4, 104. <https://doi.org/10.12688/wellcomeopenres.15341.1>

Digital Object Identifier (DOI):

[10.12688/wellcomeopenres.15341.1](https://doi.org/10.12688/wellcomeopenres.15341.1)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Publisher's PDF, also known as Version of record

Published In:

Wellcome Open Research

Publisher Rights Statement:

© 2019 Zielinski T et al. This is an open access article distributed under the terms of the Creative Commons Attribution Licence, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.





RESEARCH NOTE

The grant is dead, long live the data - migration as a pragmatic exit strategy for research data preservation [version 1; peer review: 1 approved]

Tomasz Zielinski ¹, Johnny Hay ², Andrew J. Millar ¹

¹SynthSys and School of Biological Sciences, University of Edinburgh, Edinburgh, EH9 3BF, UK

²EPCC, University of Edinburgh, Edinburgh, EH9 3FD, UK

v1 First published: 02 Jul 2019, 4:104 (<https://doi.org/10.12688/wellcomeopenres.15341.1>)
Latest published: 02 Jul 2019, 4:104 (<https://doi.org/10.12688/wellcomeopenres.15341.1>)

Abstract

Open research, data sharing and data re-use have become a priority for publicly- and charity-funded research. Efficient data management naturally requires computational resources that assist in data description, preservation and discovery. While it is possible to fund development of data management systems, currently it is more difficult to sustain data resources beyond the original grants. That puts the safety of the data at risk and undermines the very purpose of data gathering.

PlaSMo stands for 'Plant Systems-biology Modelling' and the PlaSMo model repository was envisioned by the plant systems biology community in 2005 with the initial funding lasting till 2010. We addressed the sustainability of the PlaSMo repository and assured preservation of these data by implementing an exit strategy. For our exit strategy we migrated data to an alternative public repository of secured funding. We describe details of our decision process and aspects of the implementation. Our experience may serve as an example for other projects in similar situation. We share our reflections on sustainability of biological data management and the future outcomes of its funding. We expect it to be a useful input for funding bodies.

Keywords

Data sharing, research data management, sustainable data infrastructure, exit strategy, research funding

Open Peer Review

Reviewer Status

Invited Reviewers

1

version 1

published
02 Jul 2019

report

1 **Helen Ougham**, Aberystwyth University,
Aberystwyth, UK

Any reports and responses or comments on the article can be found at the end of the article.

Corresponding author: Andrew J. Millar (andrew.millar@ed.ac.uk)

Author roles: **Zielinski T:** Conceptualization, Software, Supervision, Writing – Original Draft Preparation; **Hay J:** Software, Writing – Original Draft Preparation; **Millar AJ:** Conceptualization, Writing – Original Draft Preparation, Writing – Review & Editing

Competing interests: No competing interests were disclosed.

Grant information: This work was funded by the Wellcome Trust through a Wellcome Institutional Strategic Support Fund (ISSF3) [204804]. This work was also supported by the Biotechnology and Biological Sciences Research Council (BBSRC) through the UK Centre for Mammalian Synthetic Biology [BB/M018040].

The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Copyright: © 2019 Zielinski T *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

How to cite this article: Zielinski T, Hay J and Millar AJ. **The grant is dead, long live the data - migration as a pragmatic exit strategy for research data preservation [version 1; peer review: 1 approved]** Wellcome Open Research 2019, 4:104 (<https://doi.org/10.12688/wellcomeopenres.15341.1>)

First published: 02 Jul 2019, 4:104 (<https://doi.org/10.12688/wellcomeopenres.15341.1>)

Introduction

Open research, data sharing and data re-use have become a priority for publicly- and charity-funded research, as expressed for example in the UK Concordat on Open Research¹. Data re-use depends on reliable metadata: a detailed description of the experimental conditions, materials used, handling procedures and analysis methods. Data management goes beyond the safe storage of data, because metadata acquisition and data discovery are equally important aspects for effective digital preservation²⁻⁴. This creates a need for computational resources that can deliver such features.

Funding bodies acknowledge that data management carries significant costs and allow budgeting for data stewardship. For larger projects this permits the development of systems suitable for a particular research domain, by supporting specific data models or streamlining metadata collection. This occasionally results in the formation of a small community resource, sometimes described as “boutique repository”. Unfortunately, while it is possible to fund data infrastructure for a project, currently, there are few schemes that could support a resource beyond the timeline of the initial grant⁵. The common approach is to cover maintenance costs by “tunnelling” funds from related projects. That is not a sustainable model and puts at risk the very data that the original grant paid to preserve.

The increasing demand for data archiving induced the creation of general repositories (e.g. Figshare⁶, Zenodo⁷, Dryad⁸) and also institutional repositories (e.g. University of Edinburgh DataShare⁹, UK Data Archive¹⁰). They may lack flexibility to support all the various needs of an active project, but they are valid alternatives for data preservation. We propose to address the sustainability problem and mitigate the risk to *boutique* data by implementing an exit strategy in the form of data migration to a larger, public repository with secured funding.

Problem description

PlaSMo stands for ‘Plant Systems-biology Modelling’ and the PlaSMo portal (plasmo.ed.ac.uk) was envisioned by the plant systems biology community during a BBSRC and GARNet workshop in July 2005. The initial 2 years development was funded as part of BBSRC’s Bioinformatics and Biological Resources call in 2008 and then supported by the European Commission’s FP7 Collaborative Project TiMet (2010–2015).

The PlaSMo portal (henceforth referred to simply as ‘PlaSMo’) became a central resource for diverse plant models: general crop models, organ-level models or complex multi-component plant system models. At the time of its creation it was a unique resource for managing and sharing plant models, many of which were refactored into common, declarative languages (SBML or SimileXML). PlaSMo repository contained over 100 described models and nearly 400 data and model files.

The main features of PlaSMo were:

- Support for multiple XML model formats: SimileXMLv3, SBML Level 2 versions 1-4, Cytoscape XGMMLv1, SBGN-MLv1

- Validation of the model format
- Managing multiple versions of the model
- Each version could have its own assets: definition file, supporting data, graphical representation, bibliography, description and comments
- Public, private or group access
- Free text search

The PlaSMo portal was implemented as a typical Java web application of its time: Apache Struts 2 as the Model-View-Controller (MVC) framework with Java Server Pages (JSPs) deployed on Apache Tomcat. The choice of Java as the language and technology stack proved to be robust and convenient. For example, the backend database was migrated from DB2 to MySQL, new model formats were added, and the Struts framework major version upgraded, all as ad-hoc tasks without the original developer present. Such tasks benefited from Java features such as strong typing, rich exception handling, well-defined JAR dependencies and IDE support.

Nevertheless, there were costs in providing such a public service including system administration, software development for occasional updates and user support. The Struts MVC framework had to be updated in a timely manner due to security concerns. There were critical vulnerabilities discovered in Struts that could permit arbitrary code execution and we observed attempts to exploit them just 8 hours after their disclosure. After the funded interval, all this work was performed as an in-kind contribution.

We noticed that PlaSMo had not been attracting new users. Its user interface was outdated, and the researchers had gained other facilities for sharing, like wikis or general repositories. It seemed that the value of the PlaSMo project was in its data rather than in its portal, hence, we decided to migrate PlaSMo content into an alternative repository.

Decision process

We planned to use a repository designed specifically for biological data instead of a general one like Figshare, Zenodo or University of Edinburgh DataShare. The general resources have no features relevant to biological data (e.g. model types, difference between model and data), they also tend to have a “flat” organization structure built around a concept of datasets. We wanted to preserve the “community” aspect of PlaSMo by having its resources grouped together and available for exploratory browsing. There is a dedicated repository for biological models: BioModels^{11,12} but it accepts only public (usually published) models in SBML format, whereas PlaSMo supported earlier stages of private model development and sharing among collaborators.

We chose FAIRDOMHub as the resource to host PlaSMo data¹³. FAIRDOMHub is powered by the SEEK platform for managing systems biology data. SEEK software was developed as part of the SysMO project, a 6-year trans-European initiative

of over 100 biological research groups¹⁴. We had previously evaluated SEEK from the perspective of handling plant models, so we knew that SEEK's features aligned well with PlaSMo capabilities¹⁵. SEEK organizes assets following the Investigation, Study, Assay (ISA) structure¹⁶, offering user friendly navigation over the ISA tree. We could preserve the PlaSMo identity, utilizing the additional concept Project. Below, we refer to FAIRDOMHub when we discuss the public web data repository and to SEEK when we discuss the underlying software platform and its concepts.

We represented PlaSMo records as SEEK entities in the following way:

- Each version of PlaSMo model is represented as a separate SEEK Modelling Assay
- PlaSMo model file becomes SEEK Model
- PlaSMo images and data files become SEEK DataFiles
- SEEK Model and DataFiles are linked to a corresponding Modelling Assay
- Metadata which is not easily represented in SEEK (e.g. comments) are appended to the description text of the Modelling Assay
- For each PlaSMo model a SEEK Study is created, and the Modelling Assays representing different versions of the model are linked to that Study
- For each user who deposited a model, a SEEK Investigation is created in their name, and all Studies representing their models are linked to that Investigation (see below)
- A SEEK project named "PlaSMo Model Repository" is created and all the Investigations, Studies, Assays, Models and DataFiles are linked to it
- All SEEK entities generated for public PlaSMo models are visible to anyone in SEEK
- For private PlaSMo models the descriptions of SEEK Studies and Assays are visible to anyone in SEEK but the actual content of Model and DataFiles is hidden

The main difficulty was how to handle permissions and ownership. SEEK has a very rich and flexible access control model (in our opinion, it is the best permission model we have seen so far) and SEEK assets can be linked to user profiles as their contribution. However, to benefit from these features we would need to have SEEK accounts for all the PlaSMo users.

We could not create matching FAIRDOMHub accounts for PlaSMo users: a) we were not entitled to perform such actions on behalf of the users, b) some users already had FAIRDOMHub accounts to which they would want their assets linked. To avoid contacting all the users with a request to create FAIRDOMHub accounts, we assumed a simplified approach.

Firstly, the creator of a PlaSMo model is recorded as a text label: "other contributors" in FAIRDOMHub. Secondly, for each PlaSMo user a SEEK Investigation is created with a title matching their name. The SEEK Studies representing PlaSMo models created by a user are linked to their Investigation. In that way the models of a particular PlaSMo user can be easily accessed by navigating to the SEEK Investigation named after them in FAIRDOMHub. It also solved the issue that SEEK requires a parent Investigation for all assets and we could not create a sensible convention for this based solely on PlaSMo model description.

If a person would like to claim their models, they would contact us with their FAIRDOMHub account and we would link the whole Investigation/Study/Assay tree to that account and grant the user the manager role for those assets. That way, the models' creators could later manage their records using the SEEK UI.

PlaSMo users were always encouraged to link to their models using PlaSMo's stable URLs. In order to preserve such links, we implemented a simple URL resolver that would redirect original PlaSMo references to the appropriate records in FAIRDOMHub.

Figure 1 shows the generalized route for implementing an exit strategy for data preservation.

Implementation

We based the migration project on the existing code for the PlaSMo portal, in order to re-use the Data Access Objects (DAOs) and Data Object Model (DOM), so we only needed to implement the new data transfer logic.

We developed a Java client for programmatic communication with the SEEK REST API. Firstly, we used the available JSON request payload examples from [SEEK REST API v1.7.0](#)¹⁷ to generate a library of SEEK DOM JavaBean classes using the [jsonschema2pojo v1.0.0](#) tool¹⁸. We performed this step manually as it was a one-off project and we did not plan to keep the SEEK client in sync with the API in case it changes. Potential future work could make use of the jsonschema2pojo tool to regenerate these SEEK DOM classes automatically in the event of an update to the API.

The migration code iterates over PlaSMo models, extracting information required to generate JavaBeans corresponding to SEEK's Investigations, Studies, Assays, Models and DataFiles entities. It then invokes the client methods to create the entities inside the SEEK instance, which serialize the JavaBean objects into JSON and submit them to the API via authenticated HTTP POST requests. During our initial tests, not all of the REST calls were consistently successful, for example sometimes we observed HTTP status 500 errors caused on the server by `SQLite3::BusyException` or `AbstractController::DoubleRenderError`. For that reason, we decided to record the API calls in a way that would allow them to be 'replayed' if needed without a risk of creating duplicate entities, always yielding a consistent ISA tree within a SEEK instance.

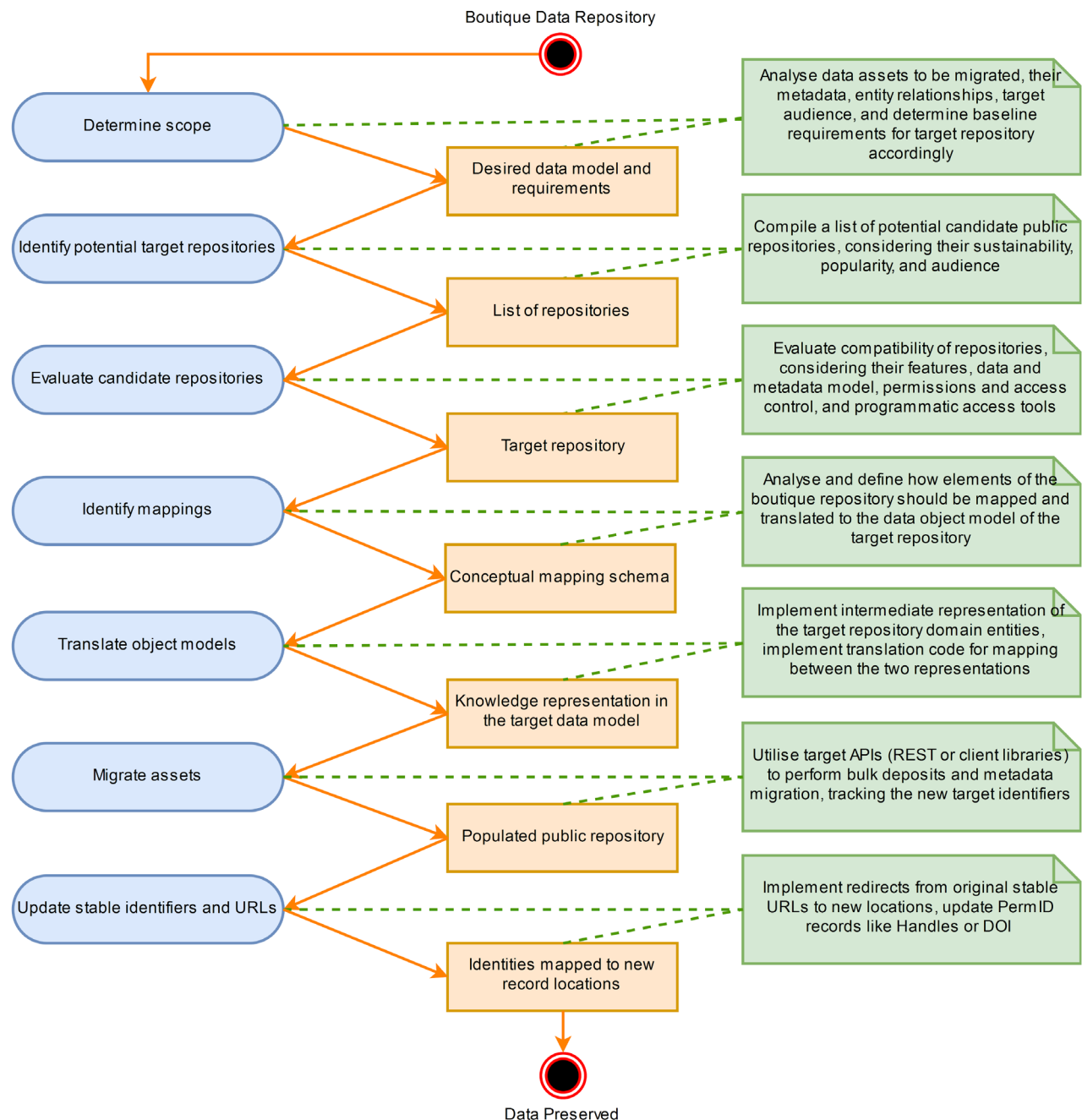


Figure 1. Implementation of an exit strategy for data preservation.

We used a local SQLite database to store the SEEK API calls, which was indexed with a GUID based on PlaSMO model UID and recorded the entire JSON payload and HTTP response for each entity in the ISA tree. This request logs database was also later used to create the mapping between PlaSMO URLs and FAIRDOMHub records (see below).

The FAIRDOMHub user interface currently does not allow for setting properties (e.g. permissions) on the whole ISA

tree, a feature necessary for our migration strategy. We implemented such bulk operations in a separate Java project, which retrieves part of the ISA tree and sets the required properties – recursively through all child entities, if desired – using the Java API client.

The recorded API calls were used to generate a mapping between PlaSMO and FAIRDOMHub identifiers. The mapping was stored as a simple csv file for ease of potential future updates.

This mapping is used by PlasmoMapper (a simple SpringBoot application) which resolves original PlaSMo URLs and redirects to the correct records in FAIRDOMHub.

Results

We performed the migration on 10th of January 2019. All the information from the PlaSMo portal are available under the [PlaSMo project on the FAIRDOMHub](#). The migration process was smooth and we did not experience any problem with the API calls. It seems that the SEEK instance on the FAIRDOMHub production server is very robust and it handles all the requests flawlessly, unlike the test SEEK's Docker containers we used during development.

Figure 2 shows the FAIRDOMHub record for version 3 of PlaSMo model 64 (Arabidopsis_clock_P2011) represented as SEEK Assay 840. The description contains the merged version specific information with the details from the main PlaSMo record (Figure 2 A, B). The other versions of that model are stored as the sibling Assays 838-841 (Figure 2 C). Each version has the list of related model and data files (Figure 2 C, D). All the model artefacts have been linked to the model owner profile in FAIRDOMHub (Figure 2 F) by running the developed SEEK Bulk Update after the migration process.

All possible PlaSMo URLs are being redirected to the corresponding records in FAIRDOMHub, for example, the main

The screenshot displays the FAIRDOMHub interface for the record 'Arabidopsis_clock_P2011 - PLM_64, version 3'. The page is divided into several sections:

- Header:** FAIRDOM HUB logo, navigation links (Browse, Help), a search bar, and user options (Register, Log in).
- Breadcrumb:** Home / Assays Index / Arabidopsis_clock_P2011 - PLM_64, version 3.
- Title:** Arabidopsis_clock_P2011 - PLM_64, version 3.
- Description (A):** This model is based on P2011 and derives from the article: **The clock gene circuit in Arabidopsis includes a repressor with additional feedback loops.** Alexandra Pokhilko, Aurora Piñas Fernández, Kieron D Edwards, Megan M Southern, Karen J Halliday & Andrew J Millar *Mol. Syst. Biol.* 2012; 8: 574, submitted to P2011 and published 6 March 2012. Link to Supplementary Information, including equations. Minor errors in the published Supplementary Information are described in a file attached to version 1 of this model (the published SBML is correct). The model describes the circuit depicted in Fig. 1 of the paper (GIF attached). It updates the Pokhilko et al. 2010 model (termed P2010), PLM_6, by including:
 - the Evening Complex genes (ELF4, ELF3, LUX),
 - light-regulated degradation of ELF3 by COP1,
 - TOC1 as a repressor rather than an activator of LHY/CCA1.
 These changes allowed the removal of hypothetical components TOC1mod (or X) and Y from the earlier models. They also reveal that the central loop of the model is a triple-repressor ring oscillator, or 'repressor' (illustrated in Fig. 8, GIF attached). SBML curation notes (please see Comments for each version).
- Contributor and Creators (F):** A section showing the contributor and creators, including a profile picture of Andrew Millar.
- Related Publications:** Alexandra Pokhilko, Aurora Piñas Fernández, Kieron D Edwards, Megan M Southern, Arabidopsis includes a repressor with additional feedback loops.. *Mol. Syst. Biol.* 8: /synopsis/msb20126.html
- Comments:** Version 2 of the model includes the StepFunction. Note that we have not yet removed functional.
- Version Comments (B):** The SBML for this version includes the lightfunction (as in Mol Syst Biol 2012 paper) to re which is used in the model simulations. However, the internal parameters of lightfunction have default Copasi ID's, parameter (photoperiod, etc.). A COPASI file for this model is also attached.
- Originally submitted to PLAStMo on 2012-03-07 12:15:12**
- Class:** Modelling Analysis
- Contributor:** BioData SynSys
- Projects:** Millar group, PlaSMo model repository
- Investigation:** Millar, Andrew (ex-PlaSMo models)
- Study:** Arabidopsis_clock_P2011 - PLM_64
- Biological problem addressed:** Gene Regulatory Network
- Organisms:** Arabidopsis thaliana
- Models:**
 - Arabidopsis_clock_P2011 - PLM_64, version 3, SUBMITTED
 - Arabidopsis_clock_P2011 - PLM_64, version 3, SIMPLIFIED
- Data (C):**
 - Fig. 1, outline of P2011 model, with P2010 inset, PLM_64_3
 - Fig. 8, cartoon illustrating repressor structure, PLM_64_3
 - ODE file for Matlab version of P2011 clock model, PLM_64_3
 - Matlab version of P2011 clock model, PLM_64_3
 - Copasi file corresponding to PLM_64 version 3, PLM_64_3
- Navigation Tree (D, E):** A tree view showing the hierarchy of models and data files. The selected item is 'Fig. 8, cartoon illustrating repressor structure, PLM_64_3 (Data file)'. The tree includes:
 - Arabis_clock_P2011 - PLM_64
 - Arabis_clock_P2011 - PLM_64, version 1
 - Arabis_clock_P2011 - PLM_64, version 2
 - Arabis_clock_P2011 - PLM_64, version 3** (selected)
 - Fig. 1, outline of P2011 model, with P2010 inset, PLM_64_3
 - Fig. 8, cartoon illustrating repressor structure, PLM_64_3
 - ODE file for Matlab version of P2011 clock model, PLM_64_3
 - Matlab version of P2011 clock model, PLM_64_3
 - Copasi file corresponding to PLM_64 version 3, PLM_64_3
 - Arabis_clock_P2011 - PLM_64, version 3, SUBMITTED
 - Arabis_clock_P2011 - PLM_64, version 3, SIMPLIFIED
 - Arabis_clock_P2011 - PLM_64, version 4
 - Arabis_clock_P2012 - PLM_70
 - At_Pokh2011v6_plasmo_ttdParams.xml - PLM_68
 - At_Pokh2011v6_plasmo_ttdParams.xml - PLM_68, version 1
 - Hidden item
 - Hidden item
 - Hidden item
 - AuxSim full - PLM_30
 - Domilant ATClock2011 - PLM_50

Figure 2. Screenshot (edited) presenting the FAIRDOMHub record for version 3 of PlaSMo model 64. A) Description of model 64 created from the main PlaSMo metadata; B) Version 3 specific details; C) List of linked model and data files; D) the navigation tree for models, versions and their data files; E) for private models the data and model files are hidden but the main metadata record is visible; F) link to the FAIRDOMHub user profile of the owner of the original PlaSMo model.

link to the model 67: http://plasmo.ed.ac.uk/plasmo/models/model.shtml?accession=PLM_64 is redirected to Study 494; the version 3 of the model: http://plasmo.ed.ac.uk/plasmo/models/model.shtml?accession=PLM_64&version=3 to Assay 840 and the file containing Matlab version of this model: http://www.plasmo.ed.ac.uk/portal_data/data/PLM64/data/98mod_P2011.m to DataFile 2499.

We believe that the plant systems biology community will benefit from the PlaSMo models migration. The models are readily available for discovery by the larger userbase of FAIRDOMHub and models can be linked to experimental data. The potential for discovery is additionally enhanced by visibility of all the descriptions even of the private models, though for private models, the actual files are not accessible. That paves the way to potential collaborations without compromising the confidentiality of the data and is only possible due to SEEK's rich permissions model. We note that this capability fulfills the stringent data sharing guidelines of UKRI-EP SRC.

We also feel that the migration boosted the profile of FAIRDOMHub as a community resource for data management and sharing. As visible in Table 1, transfer of the PlaSMo models substantially increased the number of available modelling assets (75% increase in model files). The effect of scale is an important aspect for attracting new users and the inclusion of plant models may popularise FAIRDOMHub among modellers.

Discussion: Sustainability of Biological Data Management

We imported all the PlaSMo assets into a larger community resource: FAIRDOMHub. The migration process became feasible only after SEEK's developers released the write API in 2018. This illustrates the importance of write APIs for data management systems.

The experience of shutting down a community repository, while preserving its data, challenges some of the popular views of the feasibility of Research Data Management. For example, the successful migration of all PlaSMo data could suggest that there is no need for new systems for data management.

Should funders still fund new software for data management?

In short, we believe the answer is yes.

Convincing researchers to invest the effort necessary to describe and deposit their data into a repository is the most difficult aspect of data management and a limiting factor in the wider adoption of data sharing. Data sharing can be achieved by using either “a stick” or “a carrot” approach.

The most successful “sticks” are the enforced, strict publication policies, as illustrated by the domain-specific requirements to deposit protein structural data (as in Protein Data Bank^{19,20}), sequencing data (as in GeneBank^{21,22}) or transcriptomics data (as in ArrayExpress^{23,24}). However, these repositories handle only narrow or single data types; there is consensus within each community on the minimal information criteria; in our opinion, these are the “easy” cases. For example, pdb files are in practice self-contained with metadata principally generated by equipment or processing software and require minimal interference from a scientist. Or, the deposited file represents all the results of a large, expensive experiment (e.g. microarray), so the effort in its preparation is small relative to the total effort in the experiment.

The current incentives (“carrots”) for data sharing are weak, considerably delayed in time and often accrue more to group leaders than to contracted researchers, hence they do not encourage adoption²⁵. An alternative approach is to incorporate data management into the daily research workflow, by providing immediate value to data producers in the form of increased productivity, specialized processing, visualisation or data aggregation. For example, the BioDare repository is widely used within the circadian community, but researchers primarily use the resource to access specialist software tools to analyse and visualise their timeseries data, so the fact that datasets are simultaneously deposited in the public domain is in reality a side effect of their normal work^{26,27}. This level of tool customization and integration is project/domain specific and not possible with general repositories. Consequently, we expect such “carrots” to be rare among institutional repositories that cater for many research domains.

Table 1. The total counts for each ISA entity type as they were in FAIRDOMHub before and after the PLaSMo models migration.

ISA Entity	Pre-Migration Total	Post-Migration Total	Total from PLaSMo	Percentage Increase
Investigation	197	212	15	7.6
Study	383	470	87	22.7
Assay	618	738	120	19.4
Model	255	446	191	75.0
Data File	1908	2097	189	9.9

User friendliness is the most important characteristic for successful data management. The development of user friendly solutions that facilitate research (providing the specialist “carrots” we describe above) remains a valid case for funding.

It is worth noting, that data management solutions may not need to be built entirely from scratch. One could leverage features of existing products (like for example SEEK or OpenBIS²⁸) and create plugins or integrate with them. Which approach is most cost effective and productive must be evaluated case by case, depending on the available know-how and expected user experience.

A positive example of promoting data management is the Wellcome Trust “Research Enrichment – Open Research” scheme²⁹, which funds small, add-on projects for existing grant holders to enable open research and data sharing. By presenting this as add-on funding, the implementation of data sharing is perceived as an additional opportunity, rather than in competition with core scientific activities for funding.

Can research data repositories be self-sustaining?

In majority of the cases, no^{30,31}.

The idea that domain-specific resources could often be maintained from subscription fees is unrealistic:

1. There is a problem of scale. If we advocate for resources that address particular needs of scientific projects, the underlying user base or even the entire research community may be too small to sustain a public repository financially. Conversely, repositories catering to a diverse community may gain sustainability but lack user uptake.
2. Data producers already commit their time and make a substantial effort to prepare data for deposit, so we cannot expect them to be charged for deposit on top of the work they do to contribute their data.
3. Charging for access to data is against the spirit of open research and data re-use. Funding agencies generally require public release of the results, so such a model would be an infringement of their policies.
4. Micropayment models, with small fees for extra features that one might use (e.g. minted DOI or a longer embargo period) could be acceptable to the users but it is impractical in the academic world. The research groups do not have access to credit cards to perform small payments, moreover, invoicing and accounting for such operations would be problematic and not cost-effective.

While it is possible to secure funding for a new project, there are currently only few funding schemes to maintain existing data resources. Incremental improvements to existing resources are also problematic as they typically do not pass the novelty and impact criteria.

Funding agencies should realise that maintaining existing resources may be as important as funding new science as it is the only way to enable data re-use. At the same time, the data repositories should gather metrics in order to demonstrate their value, for example numbers of active users, visits, downloads.

How to deliver data longevity?

Our PlaSMo migration demonstrates that data longevity can be achieved by implementing an effective exit strategy. In our case, we found a close match for our metadata model in FAIR-DOMHub. If not a perfect match, it is always possible to find a generic destination that can at least archive all the data. The actual implementation of a migration involves additional costs, but, in the long term, it is usually cheaper than maintaining a running resource.

The biggest value of data repositories lies in their data; hence, we would recommend creation of funding opportunities that could be used to “rescue the data”. Maybe data migration could constitute part of the income agreement for maintaining destination repositories. For example, a repository could receive extended funding on the condition that it would implement adoption of data from other projects.

Currently, data migration seems to be an inevitable reality of data preservation. Permanent identifiers (like DOIs or handles) which can resolve to the actual location facilitate this process. If PlaSMo models had DOIs we would not need to deploy PlasmoMapper to handle original URLs. Unfortunately, participation in permanent identifier schemes incurs additional financial costs, which paradoxically may accelerate the demise of a repository.

In Horizon2020, the EU funded various initiatives to provide European Research e-Infrastructure, and participating consortiums offer permanent identifiers as part of their services. Sadly, although the initiative is already centrally funded, the identifiers (handles) are provided only as a paying service. We believe that permanent identifiers should be available free of charge not only for data projects but even for individuals as a public service, similar to street address systems.

Conclusions

We shared our experience in securing the PlaSMo project’s legacy and assuring data longevity by successfully implementing an exit strategy in the form of data migration. We believe that further progress in open research and data sharing can only be achieved by integration between different resources that together can be incorporated into research workflows. We are concerned over the existing funding opportunities for data management and how they might put at risk the safety of scientific data.

Reuse potential

The Java Client for SEEK REST API and the bulk property setter, described here, can be of value for other projects. The client can be used to integrate other Java projects with SEEK,

for example to automate data deposition. The bulk property setter compensates for the currently missing feature in SEEK UI. Running the setter is currently the easiest way to publish multiple datasets constituting a research outcome. For these reasons we made the relevant code available as two separate packages.

Data availability

Underlying data

All data underlying the results are available as part of the article and no additional source data are required.

Software availability

Java Client for SEEK API

Source code is available from: <https://github.com/SynthSys/Seek-Java-RESTClient>

Archived source code at time of publication: <https://doi.org/10.5281/zenodo.3250951>³²

Licence: MIT

SEEK Bulk Update

Source code is available from: <https://github.com/SynthSys/Seek-Bulk-Update>

Archived source code at time of publication: <https://doi.org/10.5281/zenodo.3250959>³³

Licence: MIT

PlaSMo portal

Source code is available from: <https://github.com/SynthSys/PlasmoPortal>

Archived source code at time of publication: <https://doi.org/10.5281/zenodo.3250855>³⁴

Licence: MIT

Grant information

This work was funded by the Wellcome Trust through a Wellcome Institutional Strategic Support Fund (ISSF3) [204804].

This work was also supported by the Biotechnology and Biological Sciences Research Council (BBSRC) through the UK Centre for Mammalian Synthetic Biology [BB/M018040].

The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

References

- HEFCE, RCUK, UUK, *et al.*: **Concordat On Open Research Data**. 2016.
[Reference Source](#)
- Wilkinson MD, Dumontier M, Aalbersberg IJ, *et al.*: **The FAIR Guiding Principles for scientific data management and stewardship**. *Sci Data*. 2016; 3: 160018.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Wittig U, Rey M, Weidemann A, *et al.*: **Data management and data enrichment for systems biology projects**. *J Biotechnol*. 2017; 261: 229–37.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Stuart D, Baynes G, Hrynaskiewicz I, *et al.*: **Practical Challenges for Researchers in Data Sharing**. *Whitepaper*. 2018; 30.
[Publisher Full Text](#)
- Knowledge Exchange Research Data, Expert Group and Science Europe Working Group, on Research Data: **Funding research data management and related infrastructures**. 2016.
[Reference Source](#)
- Figshare.
[Reference Source](#)
- Zenodo.
[Reference Source](#)
- Dryad.
[Reference Source](#)
- Edinburgh DataShare.
[Reference Source](#)
- UK Data Archive.
[Reference Source](#)
- BioModels.
[Reference Source](#)
- Glont M, Nguyen TVN, Graesslin M, *et al.*: **BioModels: expanding horizons to include more modelling approaches and formats**. *Nucleic Acids Res*. 2018; 46(D1): D1248–53.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Wolstencroft K, Krebs O, Snoep JL, *et al.*: **FAIRDOMHub: a repository and collaboration environment for sharing systems biology research**. *Nucleic Acids Res*. 2017; 45(D1): D404–D407.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Wolstencroft K, Owen S, Krebs O, *et al.*: **SEEK: a systems biology data and model management platform**. *BMC Syst Biol*. 2015; 9(1): 33.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Troup E, Clark I, Swain P, *et al.*: **Practical evaluation of SEEK and OpenBIS for biological data management in SynthSys; first report**. University of Edinburgh. 2015.
[Reference Source](#)
- Rocca-Serra P, Brandizi M, Maguire E, *et al.*: **ISA software suite: supporting standards-compliant experimental annotation and enabling curation at the community level**. *Bioinformatics*. 2010; 26(18): 2354–6.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- SEEK REST API.
[Reference Source](#)
- Littlejohn J: **jonschema2pojo** [Internet].
[Reference Source](#)
- Protein Data Bank.
[Reference Source](#)
- Berman HM, Battistuz T, Bhat TN, *et al.*: **The Protein Data Bank**. *Acta Crystallogr Sect D Biol Crystallogr*. 2002; 58(Pt 6 No 1): 899–907.
[PubMed Abstract](#) | [Publisher Full Text](#)
- GeneBank.
[Reference Source](#)
- Benson DA, Cavanaugh M, Clark K, *et al.*: **GenBank**. *Nucleic Acids Res*. 2013; 41(Database issue): D36–42.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- ArrayExpress.
[Reference Source](#)
- Brazma A, Parkinson H, Sarkans U, *et al.*: **ArrayExpress—a public repository for microarray gene expression data at the EBI**. *Nucleic Acids Res*. 2003; 31(1): 68–71.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Van den Eynden V, Knight G, Vlad A, *et al.*: **Towards Open Research: practices, experiences, barriers and opportunities** [Internet]. 2016.
[Publisher Full Text](#)
- BioDare.
[Reference Source](#)
- Zielinski T, Moore AM, Troup E, *et al.*: **Strengths and limitations of period**

- estimation methods for circadian data. *PLoS One*. 2014; 9(5): e96462.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
28. Bauch A, Adamczyk I, Buczek P, *et al.*: **openBIS: a flexible framework for managing and analyzing complex data in biology research**. *BMC Bioinformatics*. 2011; 12: 468.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
29. **Research Enrichment – Open Research**. [Online]: Wellcome Trust;.
[Reference Source](#)
30. OECD: **Business Models for Sustainable Data Repositories**. OECD Science, Technology and Innovation - Policy Papers. 2017.
[Reference Source](#)
31. RDA-WDS Interest Group on Cost Recovery for Data Centres, Dillo I, Hodson S, *et al.*: **Income Streams for Data Repositories**. 2016.
[Publisher Full Text](#)
32. Zielinski T, Hay J: **SynthSys/Seek-Java-RESTClient: Java RestClient for SEEK API 1.7.0 (Version v1.0.0)**. *Zenodo*. 2019.
<http://www.doi.org/10.5281/zenodo.3250951>
33. Zielinski T, Hay J: **SynthSys/Seek-Bulk-Update: Bulk Update For Seek API 1.7.0 (Version v.1.0.0)**. *Zenodo*. 2019.
<http://www.doi.org/10.5281/zenodo.3250959>
34. Zielinski T, Tindal C: **SynthSys/PlasmoPortal: The last working version of PlaSMo portal (Version v2.1.5)**. *Zenodo*. 2019.
<http://www.doi.org/10.5281/zenodo.3250855>

Open Peer Review

Current Peer Review Status:

Version 1

Reviewer Report 16 July 2019

<https://doi.org/10.21956/wellcomeopenres.16751.r35896>

© 2019 Ougham H. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



Helen Ougham

Institute of Biological, Environmental and Rural Sciences (IBERS), Aberystwyth University, Aberystwyth, UK

This paper represents both a useful case study in migrating a biological research resource - in this case, the plant models repository PlaSMo - to a new public repository - and a description of the serious practical, especially financial difficulties, of maintaining specialised datasets. Data sharing, though both inherently desirable and a requirement of many funders, is not a trivial exercise in cases where no major international repository/database is available; there is often a real danger that research outputs will be either lost or made available in a form unsuitable for effective re-use.

The original PlaSMo portal, which used computational features appropriate to its time, was designed as a resource for plant and agricultural researchers to access a range of plant models, some very new, but others long-established but in danger of being lost as their originators retired and, in some cases, source code became difficult to access. It was always intended that those in the rapidly-developing systems biology community should be able to capitalise on these models to assist them in extending modelling of biological processes from the cellular to the organ, whole plant and ultimately crop level. Now, some years later, the number of plant systems biology models is much greater, and it made sense for the authors to migrate the models and datasets originally available through PlaSMo to a contemporary repository already accommodating many systems biology models and in use by current systems biologists. At the same time, this would secure the classic crop models for the plant physiologists and breeder. Carrying out this migration addressed issues of potential security threats as well as reducing the overhead inherent in maintaining the PlaSMo resource in its original form.

The paper clearly sets out the rationale for the work, the steps taken, the tools used, and the form in which the original PlaSMo models and associated files are now to be found in FAIRDOM hub; the latter has grown considerably as a result, particularly with respect to the number of models available. Although certain aspects of the original PlaSMo have been lost (ability to run web-based simulations, for example), this is an inevitable consequence of the move and is unlikely to adversely affect most current users of the models.

The paper is generally very well written; there are a few sections where the English reads a little as though

it was written by a non-native speaker, but the meaning is always very clear.

A couple of minor typos: there is one instance where PlaSMo is written PlaSMO, and this should be amended for consistency; and 'GeneBank' should be 'GenBank'.

I did notice a spelling mistake on <https://github.com/SynthSys/PlasmoPortal>, where FAIRDOMhub is shown as FaridoHub!

All the URLs in the article worked correctly at the time of this review.

The table and the two figures are useful and appropriate. In Figure 1 the boxes on the right (in green) have 'folded over' corners which in some cases slightly obscure the text; this should be easy to address.

Is the work clearly and accurately presented and does it cite the current literature?

Yes

Is the study design appropriate and is the work technically sound?

Yes

Are sufficient details of methods and analysis provided to allow replication by others?

Yes

If applicable, is the statistical analysis and its interpretation appropriate?

Not applicable

Are all the source data underlying the results available to ensure full reproducibility?

Yes

Are the conclusions drawn adequately supported by the results?

Yes

Competing Interests: I was, with one of the authors (Andrew Millar) a grantholder on the BBSRC-funded project which established the original PlaSMo repository. However I retired from Aberystwyth University in 2010 and believe that I am able to give a fair and unbiased review of this paper, which has given me an interesting and useful update on progress in the area.

Reviewer Expertise: My background is in plant science (including crop science) and bioinformatics, but not in computer science. As a grantholder on the original PlaSMo project, I am able to assess the background to the work, its current value, and the form in the PlaSMo models and associated files have been migrated and made available, but I am not able fully to assess the computational infrastructure aspects.

I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.